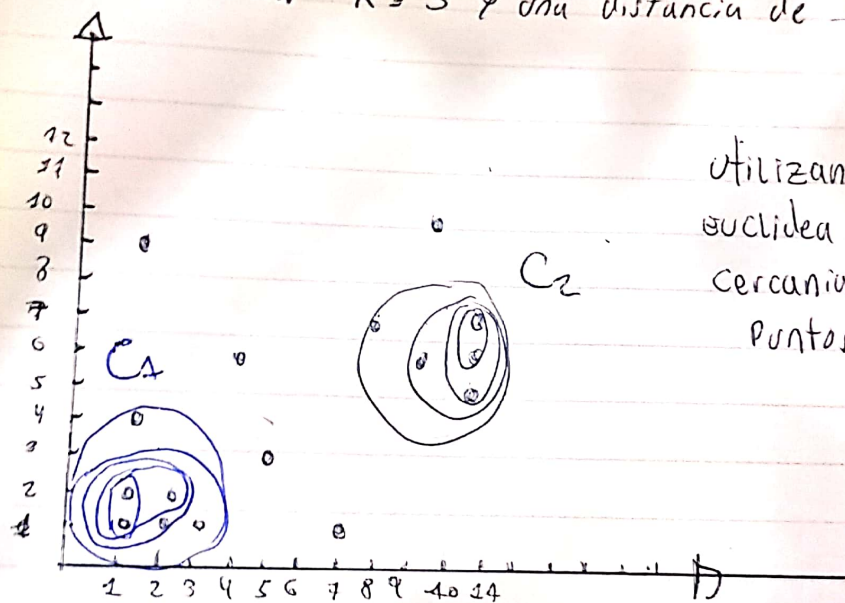


- 1) Para solucionar este problema y pudiendo visualizar los puntos, utilizaria DBSCAN para que aquellos puntos en un estado intermedio, que los considero "ruido" queden sin clasificar. ~~ahora tambien puede ser ruido ch~~

Para eso defino los hiperparametros de DBSCAN que son  $K$  y  $\epsilon$ .

Voy a tomar  $K = 3$  y una distancia de  $\epsilon = 2$



Voy a unir los <sup>puntos</sup> que están a una distancia menor a  $\epsilon = 2$  de  $K = 3$  puntos

De esta forma podria solucionarse para este caso esos puntos que no serian de ni un cluster o de otro, aunque como vemos, el punto (9, 10) deberia ir al Cluster 2 pero con la distancia que use no entra, se podria buscar una que sea mas optima.

No se si entendi bien lo que pedia el ejercicio, en caso de que esos puntos "ruidosos" tengan que ir a algun grupo si o si, usaria ... sigue...

un algoritmo como k-means tt o clustering Jerarquico  
para buscar nuevos centroides y que se ajuste a los  
puntos pedidos.



Anoto a cuantos  
nodos linka la  
P<sub>0</sub> y P<sub>1</sub>

$$\begin{matrix} P_0 = 2 & P_2 = 1 & P_4 = 2 & P_6 = 1 \\ P_1 = 1 & P_3 = 4 & P_5 = 1 & P_7 = 1 \end{matrix}$$

Hoja 3/5

prueba de  
teletr.  
a cuant

3- A) inicializando los pesos de cada nodo en  $\frac{1}{8}$ :

Lo que hago  
para cada pagina  
es, primero  
multiplicar  
A x el  
peso del nodo  
y linka a  
esa pagina,  
multiplicado  
por la probabilidad  
que le corresponde

$$\begin{aligned} P_0 &= 0.8 \times \left( \frac{1}{8} \times (0) \right) + 0.2 \cdot \frac{1}{8} = 0.025 \\ P_1 &= 0.8 \times \left( \frac{1}{8} \cdot ((1) + \frac{1}{2}) \right) + 0.2 \cdot \frac{1}{8} = \frac{7}{40} \\ P_2 &= 0.8 \times \left( \frac{1}{8} \cdot (\frac{1}{2} + \frac{1}{4}) \right) + 0.2 \cdot \frac{1}{8} = \frac{1}{10} \\ P_3 &= 0.8 \times \left( \frac{1}{8} \cdot (1) \right) + 0.2 \cdot \frac{1}{8} = \frac{1}{8} \\ P_4 &= 0.8 \times \left( \frac{1}{8} \cdot (1 + \frac{1}{4}) \right) + 0.2 \cdot \frac{1}{8} = \frac{3}{20} \\ P_5 &= 0.8 \times \left( \frac{1}{8} \cdot (\frac{1}{4} + \frac{1}{2}) \right) + 0.2 \cdot \frac{1}{8} = \frac{1}{10} \\ P_6 &= 0.8 \times \left( \frac{1}{8} \cdot (\frac{1}{4}) \right) + 0.2 \cdot \frac{1}{8} = \frac{1}{20} \\ P_7 &= 0.8 \times \left( \frac{1}{8} \cdot (\underbrace{\frac{1}{2}}_{\text{nodo 4}} + \underbrace{1}_{\text{nodo 5}} + \underbrace{1}_{\text{nodo 6}}) \right) + 0.2 \cdot \frac{1}{8} = 0.275 \end{aligned}$$

repartir teniendo  
en cuenta a las  
paginas a las  
cual le da  
peso y  
sumandose  
 $(1-\beta) \cdot \frac{1}{8}$   
que es la  
prueba de  
teletransportarse  
a cualquier  
nodo.

Vuelvo a iterar

no recibe  
de nadie

$$\begin{aligned} P_0 &= 0.8 \cdot (0) + 0.2 \cdot \frac{1}{8} = 0.025 \\ P_1 &= 0.8 \cdot \left( \frac{1}{10} \cdot \frac{1}{2} \right) + 0.2 \cdot \frac{1}{8} = 0.065 \\ P_2 &= 0.8 \cdot \left( \frac{1}{8} \cdot \frac{1}{4} + 0.025 \cdot \frac{1}{2} \right) + 0.2 \cdot \frac{1}{8} = \frac{3}{50} \\ P_3 &= 0.8 \cdot (0.275 \cdot 1) + 0.2 \cdot \frac{1}{8} = \frac{49}{200} \\ P_4 &= 0.8 \cdot \left( 1 \cdot 0.025 + \frac{1}{4} \cdot \frac{1}{8} \right) + 0.2 \cdot \frac{1}{8} = \frac{7}{100} \\ P_5 &= 0.8 \cdot \left( \frac{1}{2} \cdot \frac{3}{20} + \frac{1}{8} \cdot \frac{1}{4} \right) + 0.2 \cdot \frac{1}{8} = \frac{11}{100} \\ P_6 &= 0.8 \cdot \left( \frac{1}{8} \cdot \frac{1}{4} \right) + 0.2 \cdot \frac{1}{8} = \frac{1}{20} \\ P_7 &= 0.8 \cdot \left( \frac{1}{2} \cdot \frac{3}{20} + \frac{1}{10} \cdot 1 + 1 \cdot \frac{1}{20} \right) + 0.2 \cdot \frac{1}{8} = \frac{44}{200} \end{aligned}$$

No sigo iterando porque lo estoy haciendo muy lento

B) Detectar que B no es confiable suspenderia cambiar  
la prueba de teletransportarse (en este caso era  
 $\frac{1}{8}$ ) a 0 para el nodo B y sumarle a los demas  
nodos  $\frac{1-\beta}{171}$  con At. Paginas confiables.

Por lo que al hacer esto y el nodo Page B  
no estar linkado por ninguna otra pagina, su peso  
quedaria en 0 y nunca se llegaria a el.

4) Stream:  $\{3, 6, 6, 3, 3, 6, 5, 6\} = S$   
 $h(X) = X \bmod 32$ .

1) Hasheando cada elemento:

$$h(3) = 3 \rightarrow 00011$$

$$h(6) = 6 \rightarrow 00110$$

$$h(5) = 5 \rightarrow 00101$$

Como flajolet martin cuenta los 0, en este caso a derecha, procesando el stream resum sea r el contador de 0 a derecha

$$3 \rightarrow 00011 \quad r = 0$$

$$6 \rightarrow 00110 \quad r = 1$$

$$6 \rightarrow 00110 \quad r = 1$$

$$3 \rightarrow 00011 \quad r = 0$$

$$3 \rightarrow 00011 \quad r = 1$$

$$6 \rightarrow 00110 \quad r = 1$$

$$5 \rightarrow 00101 \quad r = 1$$

$$6 \rightarrow 00110 \quad r = 1$$

Como en FM la cont. elementos dif (orden 0) es  $2^r = 2^1 = 2$ , está mal porq el  $M^0(5)$  real es 3.

B) Una solución podría ser  $h(X) = X + 6 \bmod 32$   $a=1, b=6$

$$h(3) = 01001$$

$$h(5) = 01011$$

$$h(6) = 01100$$

Donde los contadores

quedarían en  $r=2$  al procesar

el 6 y por lo tanto  $2^2 = 4$ , que así bien sigue estimando igual de mal



Posible solución efectiva para este stream

Matías Nicolás

Hoja 5/5

Lo mejor sería colocar otra función de Hashing  
más con otro estimador para que en un momento sea  
estimado como  $z^2 = 4$  y el otro  $z^1 = 2$ . Para  
que al calcular la media de 3 y estime correctamente

c) Propongo  $a=2$  y  $b=3$

$$h(x) = 2x + 3 \text{ Mod } 32$$

$$h(3) = 01001$$

$$h(5) = 01101$$

$$h(6) = 01111$$

que como se puede ver, si

seguimos procesando el

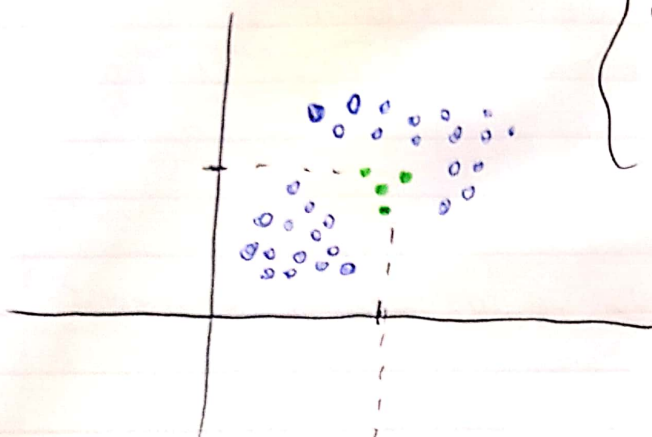
stream el contador quedará

$$\text{en } 0 \text{ y } M^0(5) = z^r = 2^0 = 1$$

No funciona

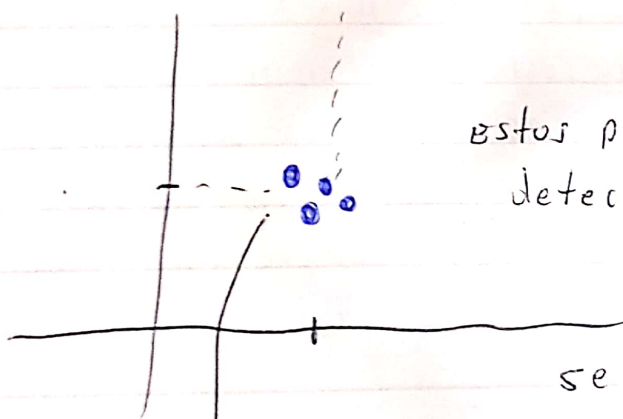
d) KNN con  $K=2$

Si tengo un modelo en el plano que entrena con los siguientes puntos:



con un set de entrenamiento desbalanceado con mayoría de puntos de clase azul.

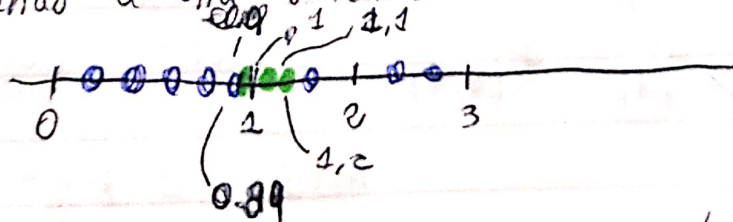
y al testearlo se utilizan los sig. puntos:



Estos puntos azules los detectaría como **verdes** debido a que el modelo está overfitteando se podría solucionar con un  $K$  más grande

↓ clasificaría mal porque los compararía con la distancia más corta de los q tiene en su entrenamiento

trascadando a una dimensión



cualquier punto que sea de clase azul pero este entre el "ruido" de 1 y 1,2 tomará los 2 puntos mas cercanos q sean 1 y 1,1 o 1,2 y lo colocará como verde.

- 2) Como necesito usar user-user para el usuario a,  
tomo los usuarios que hayan calificado la pelicula 3  
y comparo sus similitudes con a, estos usuarios  
son: c, e, f.

$$a: [1 \ 1 \ 1 \ 5 \ 1 \ 2 \ 3]$$

$$c: [2 \ 2 \ 4 \ 1 \ 2 \ 1 \ 4]$$

$$e: [1 \ 3 \ 3 \ 1 \ 1 \ 4 \ 1]$$

$$f: [4 \ 1 \ 5 \ 1 \ 1 \ 1 \ 4]$$

Pearson es restar el  
promedio de cada vector  
y Normalizarlo

Ahora tomo solo las componentes en comun con a

$$a: [1 \ 1 \ 1 \ 5 \ 1]$$

$$a = [1 \ 1 \ 5 \ 1 \ 2 \ 3] \quad , \quad \text{Prom} = 2.4$$

$$c = [2 \ 4 \ 1 \ 2 \ 4 \ 2] \quad , \quad \text{Prom} = 2.8$$

$$e = [1 \ 3 \ 1 \ 1 \ 1 \ 2] \quad , \quad \text{Prom} = 1.75$$

$$f = [4 \ 5 \ 1 \ 1 \ 4 \ 1] \quad , \quad \text{Prom} = 3$$

Resto el promedio a cada vector

$$a = [-7/5, 1, 13/5, -7/5, -2/5, 0.6]$$

$$c = [-4/5, 1.2, 1, -4/5, 1.2, -4/5]$$

$$e = [1, 1.25, 1, -0.75, -0.75, 0.25]$$

$$f = [1, 2, -2, 1, 1, -2]$$

Normalizo vectores:

$$a = [-0.418, 1, 0.7768, -0.418, -0.418, 0.179]$$

$$c = [-0.365, 0.5477, 1, -0.365, 0.5477, -0.365]$$

$$e = [1, 0.753, -0.452, -0.452, 0.15]$$

$$f = [0.267, 0.534, -0.534, 1, 0.267, -0.5345]$$



Los reescribo:

$$\begin{aligned} a &= [-0.41, 1, 0.776, -0.418, -0.11, 0.179] \\ c &= [-0.36, 0.54, 1, -0.365, 0.54, -0.365] \\ e &= [?, 0.75, 1, -0.45, -0.45, 0.15] \\ f &= [0.267, 0.53, -0.53, 1, 0.267, -0.53] \end{aligned}$$

Ahora busco las similitudes entre ellos haciendo el Producto de Vectores

$$\begin{aligned} \text{Sim}(a, c) &= 0.175 \\ \text{Sim}(a, e) &= 0.165 \\ \text{Sim}(a, f) &= -0.64 \end{aligned} \quad \left. \vphantom{\begin{aligned} \text{Sim}(a, c) &= 0.175 \\ \text{Sim}(a, e) &= 0.165 \\ \text{Sim}(a, f) &= -0.64 \end{aligned}} \right\} \begin{array}{l} \text{tomo estos 2 ya} \\ \text{que el otro es negativo} \end{array}$$

$\rightarrow$  calculo el Rating  $\left\{ \begin{array}{l} \text{valorado por user c} \\ \text{valorado por user e} \end{array} \right.$

$$=) r_{a,3} = \frac{0.175 \cdot (4) + 0.165 \cdot (3)}{0.175 + 0.165} \approx 3.51 \approx (4)$$

lo cual tiene sentido ya que la de los otros eran 4 y 3 porra era peli.



5)

b) usando 2 modelos de Regresión softmax con 2 clases (lo que sería Regresión Logística)

MA

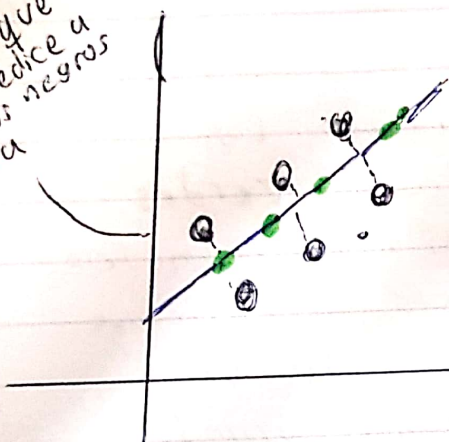
a) Regresión: usando un modelo de Regresión Lineal. Sean los

puntos para train y **test** el modelo entrenó

$$h(x) = y = mx + b$$

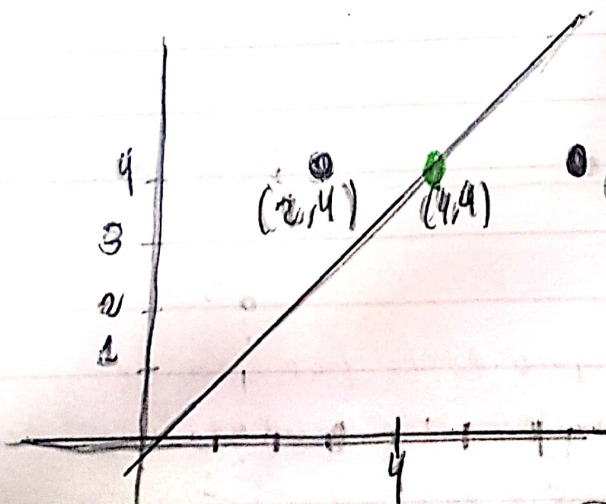
sea la recta de entrenamiento que está a igual distancia de los puntos negros de cada lado (lo que quiere decir pero deberían estar a igual distancia)

La recta que mejor predice a los puntos negros es esa



Los puntos verdes son con los cuales se entrenan, que pasan justo por la recta con lo cual no tienen error

↓ es numérico  
V



$y = x$   
Puntos: (2, 4), (6, 4)

$$MSE = \frac{1}{n} \sum (\bar{x} - x)^2$$

↓      ↓  
valor    valor  
real    que predice

Para este caso de estos 2 puntos es 0